

**APPLICATION**  
**FOR**  
**UNITED STATES LETTERS PATENT**

**APPLICANT NAME:** E. Lebrun et al

**TITLE:** METHOD AND SYSTEM FOR DETERMINING A PATH BETWEEN TWO POINTS OF AN IP NETWORK OVER WHICH DATAGRAMS ARE TRANSMITTED

**DOCKET NO.:** FR920030024US1

**INTERNATIONAL BUSINESS MACHINES CORPORATION**

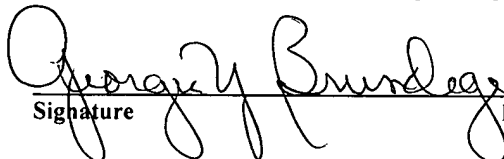
**Certificate of Mailing Under 37 CFR 1.10**

I hereby certify that, on the date shown below, this correspondence is being deposited with the United States Postal Service in an envelope addressed to the Commissioner for Patents, P.O. Box 1450, Alexandria, VA 22313-1450 as "Express Mail Post Office to Addressee"

"Express Mail" Label No.: EV 342660546 US

On: 3/25/04

Georgia Y. Brundage  
Typed or Printed Name of Person Mailing Correspondence

 3/25/04  
Signature Date

# **METHOD AND SYSTEM FOR DETERMINING A PATH BETWEEN TWO POINTS OF AN IP NETWORK OVER WHICH DATAGRAMS ARE TRANSMITTED**

## ***Background of the invention***

The present invention relates generally to computer networks, and more particularly with a method, system and computer program for determining the transmission path of datagrams sent between two points of an Internet Protocol (IP) network.

## **Internet**

The Internet is a global network of computers and computers networks which connects computers that use a variety of different operating systems or languages, including UNIX, DOS, Windows, Linux and Macintosh. To facilitate and allow the communication among these various systems and languages, the Internet uses a language referred to as TCP/IP ("Transmission Control Protocol/Internet Protocol"). TCP/IP protocol supports the following functions on the Internet:

- transmitting and receiving electronic mail,
- logging into remote computers ("Telnet"), and
- transferring files and programs from one computer to another ("FTP" or "File Transfer Protocol").

## **TCP/IP**

The TCP/IP protocol suite is named for two of the most important protocols:

- a Transmission Control Protocol (TCP), and
- an Internet Protocol (IP).

The first design goal of TCP/IP is to build an interconnection of networks that provides universal communication services: an "internetwork", or "internet". Each physical network has its own technology dependent communication interface in the form of a programming interface that provides basic communication functions running between the physical network and the user applications. The architecture of the physical networks is hidden from the user. The second design goal of TCP/IP is to interconnect different physical networks to form what appears to the user to be one large network. TCP is a transport layer protocol providing end to end data transfer. It is responsible for providing a reliable exchange of information between two computer systems. Multiple applications can be supported simultaneously over one TCP connection between two computer systems. IP is an internetwork layer protocol hiding the physical network architecture below it. Part of the task of communicating messages between computers is routing that ensures that messages will be correctly directed within the network to their destination. IP provides this routing function. An IP message is called an IP Datagram. Application Level protocols are used on top of TCP/IP to transfer user and application data from one origin computer system to one destination computer system. Such Application Level protocols are for instance File Transfer Protocol (FTP), Telnet, Gopher, Hyper Text Transfer Protocol (HTTP).

### **World Wide Web**

With the increasing size and complexity of the Internet, tools have been developed to help find information on the network, often called navigators or navigation systems. Navigation systems that have been developed include standards such as Archie, Gopher and WAIS. The World Wide Web ("WWW" or "the Web") is a recent superior navigation system. The Web is:

- an Internet-based navigation system,
- an information distribution and management system for the Internet, and
- a dynamic format for communicating on Internet.

The Web seamlessly integrates format of information, including still images, text, audio and video. A user on the Web using a graphical user interface may transparently

communicate with different host computers on the system, different system applications (including FTP and Telnet), and different information formats for files and documents including, for example, text, sound and graphics.

## **IP Router**

One of the basic functions of IP is its ability to form connections between different physical networks. This is due to the flexibility of IP to use almost any physical network below it, and to the IP routing algorithm. A system that builds connections between networks is termed a "router". A "router" is a computer that interconnects two networks and forwards messages from one network to the other. Routers are able to select the best transmission path between networks. The basic routing function is implemented in the IP layer of the TCP/IP protocol stack, so any host (or computer) or workstation running TCP/IP over more than one interface could, in theory, forward messages between networks. Because IP implements the basic routing functions, the term "IP Router" is often used. However, dedicated network hardware devices called "routers" can provide more sophisticated routing functions than the minimum functions implemented in IP.

## **IP Routing**

When data is sent to a remote destination, each IP datagram is first sent to a local router. An incoming datagram that specifies a destination IP address other than one of the local router IP address is treated as a normal outgoing datagram. This outgoing datagram is subject to the IP routing algorithm of the router, which selects the next hop for the datagram. The router forwards each datagram towards its final destination. A datagram travels from one router to another until it reaches a router connected to the destination. Each intermediate router along the end-to-end path selects the next hop used to reach the destination. The next hop represents the next router along the path to reach the destination. This next router can be located on any of the physical networks to which the intermediate router is attached. If it is a physical

network other than the one on which the host originally received the datagram, then the result is that the intermediate router has forwarded the IP datagram from one physical network to another. An "IP routing table" in each router is used to forward datagrams between networks. A basic IP routing table comprises information about the locally attached networks and the IP addresses of other routers located on these networks, plus the networks to which they attach. A routing table can be extended with information on IP networks that are farther away, and can also comprise a default route, but it still remains a table with limited information. A routing table represents only a part of the whole IP networks. A router having such a routing table is called "a router with partial routing information". A robust routing protocol must provide the ability to dynamically build and manage information in the IP routing table. As the changes in the network topology may occur, the routing tables must be updated with minimal or without manual intervention.

### **IP Path**

When a datagram is sent to a remote destination, the succession of hops that the datagram flow through, form an "IP path". The determination of the IP path is essential to understand how the datagram reaches its final destination when sent from an origin end point. Common network problems arise when the path is "broken", that means when one (or more) intermediate router along the path is not able to properly forward datagrams to the next hop. To fix such a network problem, it is necessary, first, to isolate the problem and then, to identify the first hop that does not perform any more the routing function along the path.

### **Path determination**

The problem is therefore to find a way to determine the path that a datagram takes when this datagram is transmitted between two points of an IP network. The point of the network at the origin of the path to identify, can be called "source host", and the point of the network at the end of this path, can be called "destination host". In

many cases, routing problems occur between a client who tries to access resources located on a remote server. The object of the present invention is to locate network problems between a source and a destination host, by identifying the IP path between this source and destination hosts and determining where this path is broken. Usually, in corporate networks, users do not investigate problems by themselves, and even if they do so, they would not have sufficient authority and control over the network to fix the problem. In most of the cases, users call a help desk, and the investigation of problems is done by a person named "Network Manager". This Network Manager is generally in a remote location and has enough authority to fix the network problem. The first task of the Network Manager is to isolate the problem and determine exactly where the path is broken.

### **Traceroute**

A tool commonly used today in networks to troubleshoot connectivity problems, is a software application called "traceroute". The traceroute application enables the determination of the route that datagrams follow from a source to a destination host. Traceroute is based upon UDP (User Datagram Protocol) and ICMP (Internet Control Message Protocol). Traceroute sends a series of "probe datagrams" using UDP to an IP address and waits for an ICMP reply. A first datagram with a "Time to Live" (TTL) of "one" is sent to the destination point. The first router to see the datagram decrements the TTL to zero and returns an ICMP "time exceeded" message. The datagram is discarded. In this way, the first router in the path is identified. This process is then repeated with successively larger TTL values in order to identify the series of routers in the path to the destination host. Each router through which the datagram passes similarly decrements the TTL by one, and when a router decrements the TTL to zero, it returns the ICMP message. If a datagram which is sent out does not elicit the expected ICMP reply, then the location of the break in the path can be determined by the TTL of this datagram, and the record of the other datagrams (with lower TTL values) for which the ICMP replies were received. Traceroute actually sends UDP datagrams to the destination point which reference a port number that is outside the normally used

range. This enables traceroute to determine when the destination host has been reached, that is when an ICMP "port unreachable" message is received. In other terms, UDP datagrams are sent with monotonically increasing values in the "Time To Live" (TTL) field, and the chosen UDP port is one most likely not to be in use. For each TTL value, the traceroute program sends a fixed number of datagrams (usually three), and receives the IP addresses of the routers responding. This process continues until an ICMP "port unreachable" datagram is received or some TTL threshold is reached (usually thirty). If a router receives an IP datagram and decrements the TTL to zero, then it returns an ICMP "time exceeded" message. If the IP datagram eventually reaches the destination host, an ICMP "port unreachable" message is returned. Combining the information from all the replies the traceroute program can report the whole route.

Investigating network problems is a very common task devoted to Network Managers. A typical situation is when network equipment (for example, a customer's workstation) cannot reach another network equipment (for example, a server). The end user (e.g. the customer) will call a help desk, (typically the Network Manager) in charge of investigating and resolving the problem. The problem for the Network Manager is to perform the investigation of the failure remotely; usually the Network Manager is far away from the customer's workstation. Tivoli and the Netview console are typical tools for remotely investigating a problem on a network (Tivoli and Netview are trademarks of IBM Corporation). These tools give a view of the network, and of the status of the connections between neighbor network nodes. The availability of the ports on the network nodes are continuously monitored. A failure on a port will be immediately reported on the monitoring screen, and the Network Manager will be informed of the problem. However, a port failure represents only one type of problem among all the problems that can cause a connection failure. Other problems (such as routing table problems in routers) cannot be detected using this type of tools. This is why the traceroute application is helpful for investigating on a connection failure. By determining the actual path that is used by the datagrams in the network, and by

indicating where this path is broken, traceroute allows a fast problem determination and failure isolation.

However, the usage of the traceroute application suffers from a major drawback. To execute a program such as traceroute on a computer or a system, the Network Manager needs to take control of the source host, either directly or remotely (remote logon). This implies the definition and usage by the Network Manager of passwords to take the control of the end user's workstation. The problem is that the end user (for example, the customer), may not want a third person to have access to his or her workstation for reason of confidentiality. The access by the Network Manager to a customer's workstation may be considered as an intrusion.

### ***Objects of the invention***

An object of the present invention is to enable a Network Manager to locate a break in the IP path between a source device and a destination device without taking control of the source device.

### ***Summary of the invention***

The present invention resides in a system, computer program and method for determining a transmission path of datagrams in an IP network from a source device to a destination device. The source device receives from a network manager station, a message for retrieving information related to the transmission path of datagrams to the destination device. The message comprises a destination address and a source address. The destination address is an IP address of the destination device, and the source address is an IP address of the network manager device. The source device sends a plurality of probe datagrams to a respective plurality of IP network devices along the transmission path to the destination device. Each probe datagram comprises as destination address, the IP address of the destination device and as source address, the IP address of the network manager station. Consequently, replies, if any,



to the probe datagrams are sent back directly to the network manager station by the IP network device along the transmission path, bypassing said the device. The invention also resides in the network manager station and the computer program product executing at the network manager station which sends the message to the source device, and receives and analyzes the replies to the probe datagrams to locate the break in the transmission path.

In accordance with other features of the present invention, the plurality of probe datagrams include a "time to live" (TTL) value, said value being decremented by each IP network device along the transmission path. When an IP network device decrements the TTL value to zero, the IP network device replies to the network manager station.

### ***Brief Description of the Figures***

- **Figure 1** shows a prior art technique used by Network Managers to monitor the availability of networks by testing the interfaces of network devices.
- **Figure 2** shows a prior art technique used by Network managers to determine the path between two points of a network.
- **Figure 3** shows a typical, prior art situation where a client device tries to connect to a server through a network having a failure on a connection.
- **Figure 4** shows how a prior art traceroute program determines the path from a client towards a server and locates the first failure along this path.
- **Figure 5** shows the structure of a prior art IP header layout of an IP datagram, said IP header comprising TTL and source address fields.
- **Figure 6** shows how the path between the client and the server and the position of the first network failure along the path are determined and how the information is sent back to the Network Manager according to the present invention.
- **Figure 7** is a flow chart showing the transfer of information between the different network devices (source and destination devices, routers along the IP path) and the Network Manager according to the present invention.

### ***Detailed Description of the Invention***

#### **Path determination between a client and a server from a Network Manager**

The present invention determines the transmission path of a datagram between two points in the network such as between a source host and a destination host or between a client station and a server. This determination is made from a point, such as a Network Management device or station, which is distant from the source host and client station. The present invention determines the transmission path without having to take control of the source host (client station) or manage remote logons (authorization, access lists, passwords, access control, etc...). In the following description, the source, destination and Network Management devices are named "client station", "server station" and "Network Manager station", respectively. However, these names should not be interpreted as limiting, and in other embodiments of the invention, other types of devices may be used. The method for determining the transmission path of datagrams in a network and detecting a failure on this path according to the present invention comprises the following steps :

- The Network Manager station sends a message (IP datagram) to the client station (source host) to determine the path between this client station and a server station (destination host). This message comprises the IP address of the server station.
- The client station receives the message sent by the Network Manager station and retrieves from this message the following information :
  1. the IP address of the server station (destination host);
  2. the IP address of the Network Manager station (source IP address of the incoming message);
- The client station then sends, using UDP (User Datagram Protocol), a series of "probe datagrams" to the IP address of the server station (destination host). IP datagrams carrying the UDP datagrams are sent with monotonically increasing values in the "time to live" (TTL) field of the IP header. For each TTL value, a fixed number of datagrams are sent. However instead of inserting in each outgoing datagram, its own source address (client station IP address), the client station inserts the IP address of the Network Manager station as source address. This operation is called "spoofing".

- If a router receives one of the IP datagrams sent by the client station and decrements the TTL to zero, then this router returns an ICMP (Internet Control Message Protocol) "time exceeded" datagram. The datagram returned is sent back according to the source address of the datagram. In this case, the returned datagram is directly sent back to the Network Manager station, because of the spoofing technique described in the previous paragraph.
- If an IP datagram sent by the client station, eventually reaches the server station (destination host), an ICMP "port unreachable" datagram is returned to the Network Manager station.

In the meantime the Network Manager station, combining the information from all the replies, can report the whole route up to the point where the first network failure occurs.

### **Network monitoring using Ping application**

There are two simple and widely used applications that are based on ICMP - Ping and Traceroute. Ping uses the ICMP Echo and Echo reply messages to determine whether a host is reachable. Ping is the simplest of all TCP/IP applications. It sends one or more IP datagrams to a specified destination host requesting a reply and measures the round trip time.

Figure 1 shows a prior art method based on the Ping application to monitor, from a remote Network Manager station (103), the availability of a network connection between a client station (104) (source host) and a server station (105) (destination host). This method comprises the following steps:

- From its work station, the Network Manager (Help Desk) (103) sends a probe IP datagram called "ping" (more properly named "Echo Request") towards each physical interface (port) of the network devices located on the path that the Network Manager wants to verify (101).

- Upon reception of this datagram (if the interface is operational), each network device sends back a reply (sometimes called "Ping Response" but more properly named "Echo Reply" ) (102).
- The Network Manager station waits for the network devices' replies. A reply is received when the network device is operational. In case of failure, no reply is received and the Network Manager station displays a "timeout" message after a given period of time.

By checking all individual interfaces along the path, the Network Manager can know whether or not the interfaces are working properly. However, this method has some drawbacks. In complex networks, the path taken by datagrams between a client station and a server station is not always known. Many different paths can be taken, and the selected path depends on complex IP routing mechanisms. Furthermore, the path can be chosen dynamically, if a dynamic routing protocol is used. Thus, it is sometimes difficult to correlate an interface failure with a connection problem, i.e. the interface failure can be unrelated to the problem when for instance this failure is not on the path. It is absolutely necessary to determine the path before verifying whether or not the interfaces along said path are working properly. In addition, a connection may be broken at a logical level while all physical interfaces are operational. If, for example, the routing function of one of the intermediate network devices along the path is not properly working, datagrams may be lost, although all interfaces are correctly working. This is the reason why, it is very useful to have means to determine the path that is actually used by the datagram flow.

### **Determination of a Datagram Path with Traceroute**

Figure 2 Shows a prior art method based on the traceroute application used by Network Managers to determine the path used by a datagram sent by a client to a server. The Network Manager first takes control of the client station with a remote login (101). Then, the Network Manager launches in the client station, the traceroute program. As described earlier, the traceroute program sends probe datagrams along

the path leading to the server station (209), causing each of the intermediate network devices (nodes, routers, ...) (210) to send back a reply ( (202) to (206) ). The client station (208) gathers the replies sent by the network devices, determines the entire path and transmits the information to the Network Manager (Help Desk) (207) .

### **Failure on the Path between Client and Server**

Figure 3 shows a typical, prior art situation where a client station (303) tries to reach a server station (304) without success due to a failure in the network. The path taken by the datagram is symbolized by the arrows (301). This path is interrupted at the point (302) where the network is broken. At this point, datagrams cannot be further transmitted.

### **Failure Investigation with Traceroute**

Figure 4 shows how the traceroute application operates to investigate a failure in a network according to prior art.

- Probe datagrams are sent by the client station towards the server station.
- Intermediate network devices (nodes, routers, ... ) reply to the client station, up to the network failure (401 to 403)
- The client station collects the replies and determines the network path up to the network failure.

The last datagram sent by the client station, reaches the network failure and is lost. No reply (404) is sent back to the client station. After a predefined time period, the client station presumes that a network failure is present just after the path partially identified. The information is forwarded to the Network manager which can easily and quickly locate the failure.

### **Spoofing Technique**

Figure 5 shows the structure of an IP datagram header. The fields in the IP header mean the following:

- **VERS:** This field identifies the IP protocol version.
- **HLEN:** The length of the IP header counted in 32-bit quantities. This length does not include the data field.
- **Service Type:** The service type is an indication of the quality of service (QoS) requested for the IP datagram.
- **Total Length:** The total length of the datagram, header and data.
- **Identification:** A unique number assigned by the sender to aid in reassembling a fragmented datagram. Each fragment of a datagram has the same identification number.
- **Flags:** This field comprises control flags
- **Fragment Offset:** This field is used to aid in reassembling the full datagram. The value of this field is the number of 64-bit segments (header bytes are not counted) that are contained in earlier fragments. If this is the first (or only) fragment, this field has a value of zero.
- **Time to Live (TTL):** This field specifies the maximum time period (in seconds) during which the datagram can travel. Theoretically, each router processing this datagram is supposed to subtract its processing time from this field. In practice, a router processes the datagram in less than one second. Thus the router subtracts one from the value indicated in this field. The TTL has become a hop-count metric rather than a time metric. When the value reaches zero, it is assumed that this datagram is traveling in a closed loop and is discarded. The initial value should be set by the higher level protocol that creates the datagram.
- **Protocol Number:** This field indicates the higher level protocol to which IP should deliver the data in this datagram.
- **Header Checksum:** This field is a checksum for the information contained in the header. If the header checksum does not match the contents, the datagram is discarded.
- **Source IP Address:** The 32-bit IP address of the host sending this datagram.

- **Destination IP Address:** The 32-bit IP address of the destination host for this datagram.

Among the IP header fields described here above, the source address is different in the present invention. In the prior art, this address is normally the IP address of the device that has originated the datagram. This IP address is used by a number of programs which need to know the origin of the datagrams. For example, a program which needs to send a message of error back to the originator of the datagram will use this IP address. This is the case with the error message sent back when the time to live field (TTL) reaches zero. This mechanism is used by the traceroute application, in particular, to collect the replies to its probe datagrams. However, in accordance with the present invention, the “spoofing” technique consists of changing the source IP address of the datagram to make the programs processing the datagram think that the origin of the datagram is different than what it is in reality. The error messages are sent back to the new spoofed address and not to the real originator (source host) of the datagram.

### **Method for Determining the Path and the Failure along the Path according to the Present Invention**

Figure 6 shows how the path between the client station (610) and the server station (611) and the position of the first failure (609) along the path are determined. Figure 6 also shows how the information is sent back to the Network Manager station (612) according to the present invention. The method of the present invention comprises the following steps:

- A message (601) is sent from the Network Manager station (612) by computer programming within the Network Manager station to the client station (610) to launch the process.
- Upon receipt of the message (601) sent by the Network Manager station, the client station (including computer programming executing within the client station) sends



probe datagrams toward the server station (611) with increasing TTL values similarly to the traceroute program (602 to 604). These probe datagrams (using the “spoofing” technique) include, as source address, the IP address of the Network Manager station (instead of the IP address of the client station).

- The intermediate routers (613) send their replies back (when TTL = 0) according to the traceroute program described above. However; because the intermediate routers rely on the spoofed address inserted in the probe datagrams, the replies (606 to 608) are sent back directly to the Network Manager station.
- The replies are collected by computer programming within the Network Manager station, which determines the first part of the path up to the network failure (or the complete path up to the server station if there is no network failure).
- The last probe datagram (605) sent by the client station (610) reaches the network failure (609) and is lost. The Network Manager station is then able to display the network path up to the network failure (609).

Figure 7 is a flow chart of the foregoing process showing the interaction between the various components of the network.

- The Network Manager station (708) sends a command (701) to the client station (709) to start the process and then waits for the replies (702) sent back by routers (702).
- The client station (709), upon receipt of the start command (701), sends the probe datagrams (703 705 707 etc...) with as source address, the spoofed address of the Network Manager station (708) and as destination address, the IP address of the server station.
- The intermediate routers (710) reply back directly to the Network Manager station (704, 706, etc...) when the TTL = 0.
- The replies (702) are collected by the Network Manager station.
- At the end of the process, the Network Manager station (708) displays (706) the path that has been computed based on the replies sent by the routers (710).

While the invention has been particularly shown and described with reference to a preferred embodiment, it will be understood that various changes in form and detail may be made therein without departing from the spirit, and scope of the invention.